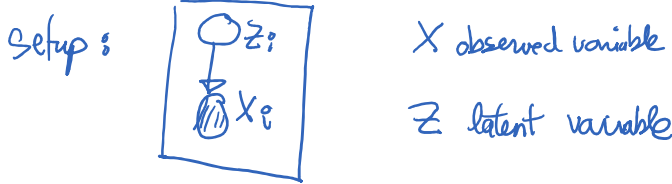


today : • EM algorithm
• GMM

EM - maximum likelihood in latent variable model



log-likelihood $\log p(x_{1:n}; \theta) = \log \left(\prod_{i=1}^n p(x_i; \theta) \right)$
 $= \sum_{i=1}^n \log p(x_i; \theta)$
 $= \sum_{i=1}^n \log \left[\sum_{z_i} p(x_i, z_i; \theta) \right]$

problem? → gives multi-modal opt. problem (non-concave)

options for ML in latent variable model

1) do gradient ascent on a non-concave obj.

f is convex
↔ -f is concave

2) EM alg. → block-coordinate ascent on auxiliary fct. which lower bounds $\log p(x_{1:n}; \theta)$
 nice interpretation in terms of filling "missing data"

i.e. E step → fill z with "soft-values"

M step → max w.r. to θ for fully observed model

trick overview:

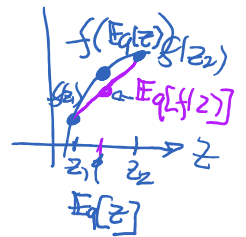
$$\log \sum_z p(x, z) = \log \sum_z q(z) \frac{p(x, z)}{q(z)}$$

$$= \log \mathbb{E}_q \left[\frac{p(x, z)}{q(z)} \right]$$

Jensen's inequality

$$\mathbb{E}_q[f(z)] \leq f(\mathbb{E}_q[z])$$

when f is concave



Jensen's inequality trick

$$\mathbb{E}_q \left[\log \frac{p(x, z)}{q(z)} \right] = \sum_z q(z) \log p(x, z) - \sum_z q(z) \log q(z)$$

$$\triangleq \mathcal{L}(q, \theta) \triangleq \underbrace{\mathbb{E}_q [\log p(x, z; \theta)]}_{\text{"expected complete log-likelihood"}} + \underbrace{H(q)}_{\text{"entropy of q"}}$$

we have $\log p(x; \theta) \geq \mathcal{L}(q, \theta) \quad \forall q, \theta$ "expected complete log-likelihood" "entropy of q"

EM algorithm: E step: $q_{t+1} \triangleq \operatorname{argmax}_q \mathcal{L}(q, \theta_t) \Rightarrow q_{t+1}(z) = p(z|x; \theta_t)$

M step: $\theta_{t+1} \triangleq \operatorname{argmax}_{\theta} \mathcal{L}(q_{t+1}, \theta)$
 $= \operatorname{argmax}_{\theta} \mathbb{E}_{q_{t+1}(z)} [\log p(x, z; \theta)]$

this is another ML problem, but for complete information

(often, replace z with $\mathbb{E}_{q_{t+1}}[z]$ in this expression)

from Jensen's ineq.

* we had $\log p(x; \theta) \geq \mathcal{L}(q, \theta)$ $\log(\mathbb{E}_q[g(z)]) \geq \mathbb{E}_q \log(g(z))$
 in Jensen's ineq., you get strict ineq. unless the dist. is degenerate (when f is strictly concave) (i.e. takes only one value)
 i.e. when $g(z) = \text{constant}$

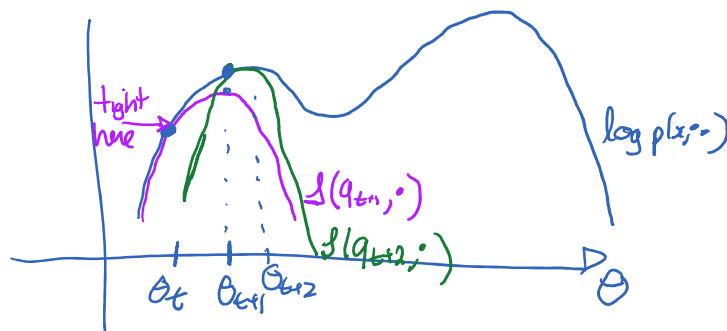
$g(z) = \text{constant}$ above i.e. $p(x, z) = \text{const.} \quad \forall z \Rightarrow q(z) \propto p(x, z)$

i.e. $q(z) = p(z|x; \theta)$ makes equality $\mathcal{L}(q, \theta) = \log p(x; \theta)$

$\mathcal{L}(q_{t+1}, \theta_t) = \log p(x; \theta_t) \geq \mathcal{L}(q, \theta_t)$ i.e. $\operatorname{argmax}_q \mathcal{L}(q, \theta_t) = p(z|x; \theta_t)$
 $\Rightarrow q_{t+1}$ maximizes $\mathcal{L}(q, \theta_t)$ w.r. to q and $\mathcal{L}(q_{t+1}, \theta_t) = \log p(x; \theta_t)$

properties:

a) $\log p(x; \theta_{t+1}) \geq \log p(x; \theta_t)$ proof: $\log p(x; \theta_{t+1}) \geq \mathcal{L}(q_{t+1}, \theta_{t+1})$



$\geq \mathcal{L}(q_{t+1}, \theta_t)$
 $= \log p(x; \theta_t)$

b) Θ_t in EM converges to a stationary pt. of $\log p(x; \Theta)$
 i.e. $\nabla_{\Theta} \log p(x; \Theta) \Big|_{\hat{\Theta}} = 0$

like k-means, initialization is crucial
 → usually do random restarts

for GMM, could use k-means++ to initialize the μ 's

14h34

c) $\mathcal{L}(q, \Theta) = \mathbb{E}_q \left[\log \frac{p(x, z; \Theta)}{q(z)} \right]$

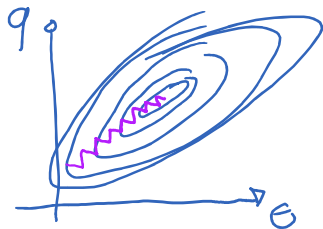
$\log p(x; \Theta) - \mathcal{L}(q, \Theta) = -\mathbb{E}_q \left[\log \frac{p(x, z; \Theta)}{q(z) p(x; \Theta)} \right]$

$\log p(x; \Theta)$
 $\mathcal{L}(q, \Theta)$
 } $KL(q \| p(\cdot | x, \Theta))$

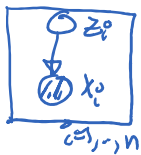
$= \mathbb{E}_q \left[\log \left(\frac{q(z)}{p(z|x, \Theta)} \right) \right] \cong KL(q(\cdot) \| p(\cdot | x, \Theta))$

KL divergence

we will revisit for variational inference $q \in \mathcal{Q}$



for GMM model



$z_i \sim \text{Mult}(\pi)$

$x_i | z_i = j \sim N(\mu_j, \Sigma_j)$

shorthand to say $z_{ij} = 1$

$\Theta = (\pi, (\mu_j)_{j=1}^k, (\Sigma_j)_{j=1}^k)$

notation: $x = x_{1:n}$

$z = z_{1:n}$

exercise:

$p(z|x) = \prod_i p(z_i|x) = \prod_i p(z_i|x_i)$

complete log-likelihood:

$\log p(x, z; \Theta) = \sum_i \left[\log p(x_i | z_i; \Theta) + \log p(z_i; \Theta) \right]$

$$\log p(x, z; \theta) = \sum_{i=1}^n \left[\log p(x_i | z_i; \theta) + \log p(z_i; \theta) \right]$$

↓ Gaussian
↓ multinomial

$$= \sum_{i=1}^n \left[\sum_{j=1}^k z_{ij} \log N(x_i | \mu_j, \Sigma_j) + \sum_{j=1}^k z_{ij} \log \pi_j \right]$$

$$\mathbb{E}_q [\log p(x, z; \theta)] = \sum_{i=1}^n \sum_{j=1}^k \mathbb{E}_q [z_{ij}] (\log N(x_i | \mu_j, \Sigma_j) + \log \pi_j)$$

$\mathbb{E}_q [z_{ij}] = q(z_{ij}=1)$ [marginal distribution]

during EM, $q_{t+1}(z) = p(z | x; \theta_t)$

weight $\gamma_{ij}^t \triangleq p(z_{ij}=1 | x_i; \theta_t) = q_{t+1}(z_{ij}=1)$

E-step is computing $q_{t+1}(z) \triangleq p(z | x; \theta_t)$

$$= \prod_i p(z_i | x_i; \theta_t)$$

$$\Rightarrow q_{t+1}(z_i) \propto p(x_i | z_i; \theta_t) p(z_i; \theta_t)$$

↓ Gaussian
↓ π_j^t

$$\gamma_{i,j}^t = q_{t+1}(z_{ij}=1) = \frac{\pi_j^{(t)} N(x_i | \mu_j^{(t)}, \Sigma_j^{(t)})}{\sum_{l=1}^k \pi_l^{(t)} N(x_i | \mu_l^{(t)}, \Sigma_l^{(t)})} \left. \begin{array}{l} \} p(x_i, z_{ij}=1 | \theta^{(t)}) \\ \} p(x_i | \theta^{(t)}) \end{array} \right\}$$

E step for GMM: compute $\gamma_{i,j}^{(k)}$ for $i=1, \dots, n$ using $\theta^{(k)}$

M step: $\max_{\{\mu_j, \Sigma_j, \pi_j\}} \sum_i \sum_j \gamma_{ij}^{(k)} [\log p(x_i | \mu_j, \Sigma_j) + \log \pi_j]$

exercise:

M step for EM for GMM

$$\hat{\pi}_j^{(k+1)} = \frac{\sum_i \gamma_{ij}^{(k)}}{n}$$

"soft count"

$$\hat{\mu}_j^{(k+1)} = \frac{\sum_{i=1}^n \gamma_{ij}^{(k)} x_i}{\sum_{i=1}^n \gamma_{ij}^{(k)}}$$

$$\hat{\Sigma}_j^{(k+1)} = \frac{\sum_{i=1}^n \gamma_{ij}^{(k)} (x_i - \hat{\mu}_j^{(k+1)}) (x_i - \hat{\mu}_j^{(k+1)})^T}{\sum_{i=1}^n \gamma_{ij}^{(k)}}$$

$$\Sigma_j^{(t+1)} = \frac{\sum_{i=1}^n \pi_{ij}^{(t)} (x_i - \mu_j^{(t+1)})(x_i - \mu_j^{(t+1)})^T}{\sum_{i=1}^n \pi_{ij}^{(t)}}$$

initialise: e.g. $\mu_j^{(0)}$ from k-means++

$\Sigma_j^{(0)}$ big spherical covariance $\Sigma_j^{(0)} = \sigma^2 I$
↑ big

$\pi_j^{(0)}$ = proportions from k-means++

EM step in GMM with fixed $\Sigma_j = \sigma^2 I$ with $\sigma^2 \rightarrow 0$

\leadsto get k-means alg?